

Regressão

Inteligência na Web e Big Data

Fabricio Olivetti de França e Thiago Ferreira Covões
folivetti@ufabc.edu.br, thiago.covoes@ufabc.edu.br

Centro de Matemática, Computação e Cognição
Universidade Federal do ABC



Regressão Linear

- Uma das formas de aprendizado de máquina é definido como um conjunto de dados na forma $\{(\mathbf{x}, y)\}$, com $\mathbf{x} \in \mathbb{R}^d$ um vetor de atributos e $y \in \mathbb{R}$ uma variável alvo, queremos descobrir um mapa entre um certo \mathbf{x} para seu valor y correspondente:

$$f : \mathbf{x} \rightarrow y.$$

- Vamos tomar como exemplo um corretor que quer aprender a avaliar o valor de um imóvel.
- Cada imóvel pode ser representado por diversas características como:
 - Metragem quadrada
 - Vagas na garagem
 - Andar
 - Bairro
 - etc.

- Para aprender a avaliar novos imóveis ele investiga pelo classificado do jornal vários exemplos:

m ²	vagas	andar	valor
80	1	2	100.000,00
120	2	5	400.000,00
100	1	10	350.000,00
...

- O objetivo é aprender:

$$f(m^2, \text{vagas}, \text{andar}) = \text{valor},$$

para quaisquer valores de m^2 , vagas e andar.

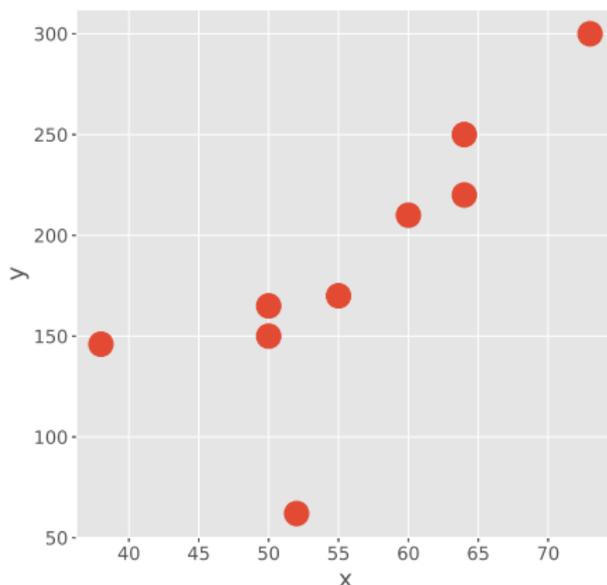
Aprendendo com exemplos

- Reduzindo nossa tabela para duas variáveis, temos:

m^2	mil \$
52	62
38	146
50	150
50	165
55	170
60	210
64	220
64	250
73	300

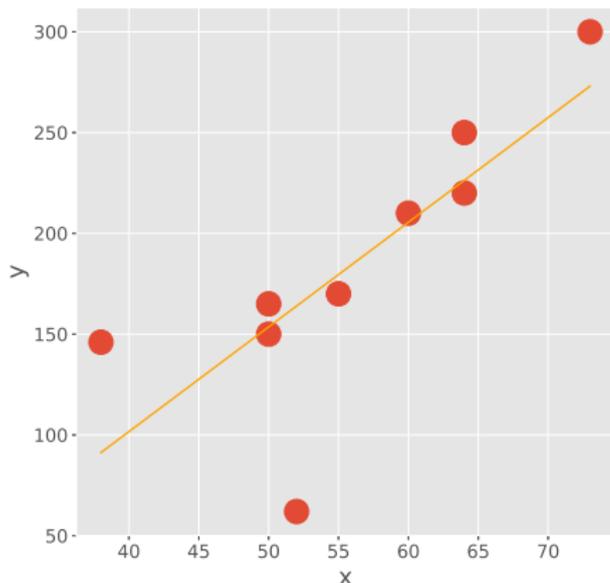
Aprendendo com exemplos

- Com duas variáveis é simples verificar a relação visualmente:



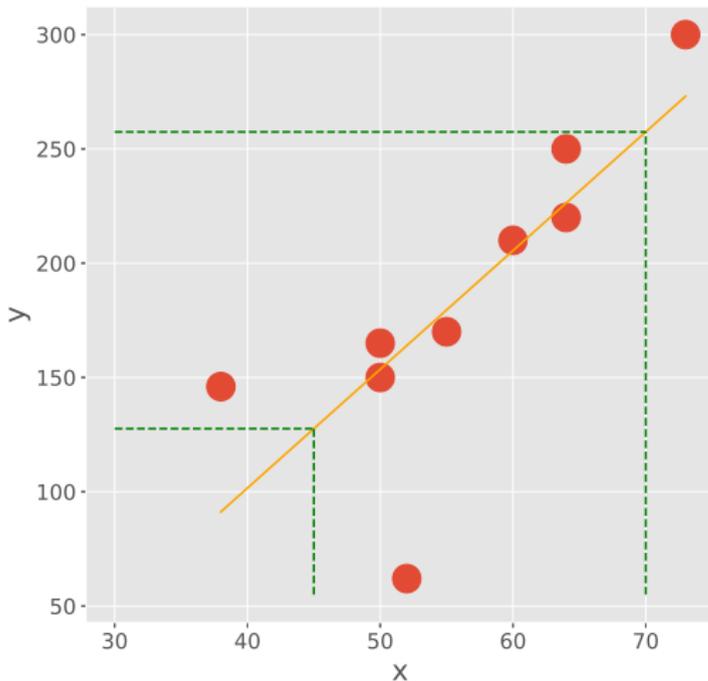
Aprendendo com exemplos

- Com duas variáveis é simples verificar a relação visualmente:



Aprendendo com exemplos

- Com essa reta podemos determinar novos valores:



- Esse método é conhecido como regressão linear.
- Ele pode ser usado quando nossas variáveis apresentam uma correlação linear entre elas!
 - mas tem alguns truques para permitir correlações não-lineares também

- Dado um conjunto com n exemplos $\{(\mathbf{x}_i, y_i)\}$, com $\mathbf{x}_i \in \mathbb{R}^d, y \in \mathbb{R}$. Queremos encontrar $f: \mathbb{R}^d \rightarrow \mathbb{R}$, tal que:

$$f(\mathbf{x}_i) \approx \mathbf{w} \cdot \mathbf{x}_i + b = y_i$$

- Renomenado nossas variáveis do exemplo, temos:

$$x_{i,1} = m^2$$

$$x_{i,2} = \text{vagas}$$

$$x_{i,3} = \text{andar}$$

$$y_i = \text{vagas}$$

$$x_i = (x_{i,1}, x_{i,2}, x_{i,3})$$

$$f(\mathbf{x}_i) = y_i$$

$$\mathbf{w} \cdot \mathbf{x}_i + b = y_i$$

$$w_1 \cdot x_{i,1} + w_2 \cdot x_{i,2} + w_3 \cdot x_{i,3} + b = y_i$$

- Definindo $w_0 = b$ e $x_{i,0} = 1$ para todo i , podemos escrever:

$$w_0 \cdot x_{i,0} + w_1 \cdot x_{i,1} + w_2 \cdot x_{i,2} + w_3 \cdot x_{i,3} = y_i$$

$$\mathbf{w} \cdot \mathbf{x}_i = y_i,$$

- ou seja, y é o produto interno entre \mathbf{w} e \mathbf{x}_i .

- Dados os n exemplos, podemos escrever na forma matricial, $X \in \mathbb{R}^{n \times (d+1)}$, $\mathbf{y} \in \mathbb{R}^{n \times 1}$:

$$X \cdot w = y,$$

- com $\mathbf{w} \in \mathbb{R}^{(d+1) \times 1}$.
- Esse modelo de regressão é dito interpretável pois é fácil determinar o quanto uma mudança em determinado x_j afeta o valor de y .

Ordinary Least Squares

Ordinary Least Squares

- O nosso problema consiste em determinar \mathbf{w} de tal forma a minimizar o erro de aproximação:

$$e(\mathbf{w}) = (y - X \cdot \mathbf{w})^T (y - X \cdot \mathbf{w}).$$

Ordinary Least Squares

- Para determinar a solução, calculamos a derivada e igualamos a zero para encontrar o mínimo local:

$$\frac{\partial e(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial (y - X \cdot \mathbf{w})^T (y - X \cdot \mathbf{w})}{\partial \mathbf{w}}.$$

- Calcule:

$$\frac{\partial e(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial (y - X \cdot \mathbf{w})^T (y - X \cdot \mathbf{w})}{\partial \mathbf{w}}.$$

$$\begin{aligned}\frac{\partial e(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial (y - X \cdot \mathbf{w})^T (y - X \cdot \mathbf{w})}{\partial \mathbf{w}} \\ \frac{\partial e(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial (yy^T - (X\mathbf{w})^T y - y^T X\mathbf{w} + (X\mathbf{w})^T X\mathbf{w})}{\partial \mathbf{w}} \\ \frac{\partial e(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial (yy^T - 2\mathbf{w}^T X^T y + \mathbf{w}^T X^T X\mathbf{w})}{\partial \mathbf{w}} \\ \frac{\partial e(\mathbf{w})}{\partial \mathbf{w}} &= -2X^T y + 2X^T X\mathbf{w}\end{aligned}$$

$$-2X^T y + 2X^T X w = 0$$

$$2X^T X w = 2X^T y$$

$$X^T X w = X^T y$$

$$w = (X^T X)^{-1} X^T y$$

- Temos então que:

$$w = (X^T X)^{-1} X^T y,$$

- com $(X^T X)^{-1} X^T$ sendo a pseudo-inversa de X .

Ordinary Least Squares

- Para dimensões grandes de X o custo da multiplicação de matrizes pode se tornar proibitivo.
- Uma implementação simples, assumindo a existência de uma função `inv` que inverte uma matriz é dada por:

```
1  ols :: [[Double]] -> [Double] -> [Double]
2  ols x y = multMtx pseudoInv (multMtx (transpose x) y)
3  where
4  pseudoInv = inv $ multMtx (transpose x) x
```

Ordinary Least Squares

- A multiplicação $X^T X$ pode ser simplificada como a soma dos produtos externos de cada x_i por ele mesmo:

```
1  ols :: [[Double]] -> [Double] -> [Double]
2  ols x y = multMtx pseudoInv (multMtx (transpose x) y)
3  where
4      pseudoInv = inv $ sum $ map autoOuterProd x
5      autoOuterProd xi = outerProd xi xi
```

Ordinary Least Squares

- Da mesma forma $X^T y$ pode ser simplificada como a soma dos vetores $x_i \cdot y_i$:

```
1  ols :: [[Double]] -> [Double] -> [Double]
2  ols x y = multMtx pseudoInv $ sumVec (y .*.. x)
3  where
4      pseudoInv = inv $ sum $ map autoOuterProd x
5      autoOuterProd xi = outerProd xi xi
```

Ordinary Least Squares

- Em um contexto de computação distribuída, imaginamos que temos pedaços de (x_i, y_i) espalhados pelo cluster de máquinas. Dessa forma podemos paralelizar pensando em funções de map e reduce:

```
1 mapper :: ([Double], Double) -> ([[Double]], [Double])
2
3 reducer :: ([[Double]], [Double]) -> ([[Double]],
4     ↪ [Double]) -> ([[Double]], [Double])
5
6 ols :: ChunksOf ([Double], Double) -> [Double]
7
8 ols xy = multMtx pseudoInv xTy
9     where
10         xTx, xTy = mapReduce mapper reducer xy
11         pseudoInv = inv xTx
```

Ordinary Least Squares

- Como implementar as funções mapper e reducer?

```
1 mapper :: ([Double], Double) -> ([[Double]], [Double])
2
3 reducer :: ([[Double]], [Double]) -> ([[Double]],
4   ↪ [Double]) -> ([[Double]], [Double])
5
6 ols :: ChunksOf ([Double], Double) -> [Double]
7 ols xy = multMtx pseudoInv xTy
8   where
9     xTx, xTy  = mapReduce mapper reducer xy
10    pseudoInv = inv xTx
```

Ordinary Least Squares

- Como implementar as funções mapper e reducer?

```
1 mapper :: ([Double], Double) -> ([[Double]], [Double])
2 mapper (x, y) = (outerProd x x, y *. x)
3
4 reducer :: ([[Double]], [Double]) -> ([[Double]],
5     ↪ [Double])
6     -> ([[Double]], [Double])
7 reducer (a, b) (a, b) = (a ..+.. b, a .+. b)
```

Gradiente Descendente

-Uma outra forma de resolver o problema é utilizando métodos de gradiente para otimização.

- Dada uma função $f(x)$ definida e diferenciável em torno de um ponto x_0 , sabemos que ela decresce mais rapidamente na direção oposta do gradiente $\nabla f(x_0)$.

- Ou seja, fazendo:

$$x^{t+1} = x^t - \alpha \cdot \nabla f(x^t).$$

- Temos que $f(x^{t+1}) \leq f(x^t)$, para um α pequeno.

- Temos então que:

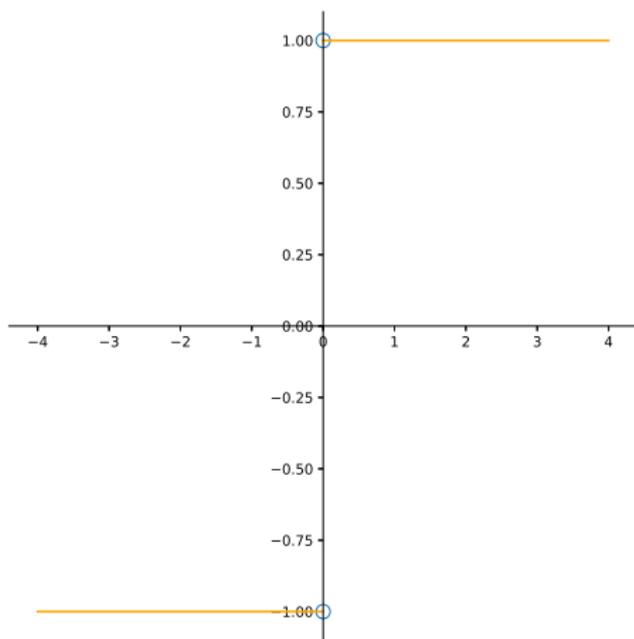
$$f(x^0) \geq f(x^1) \geq f(x^2) \geq f(x^t).$$

- Estamos interessados na minimização da aproximação de $f(\mathbf{x})$ com y :

$$e(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n |y_i - \mathbf{x}_i \cdot \mathbf{w}|,$$

Gradiente Descendente

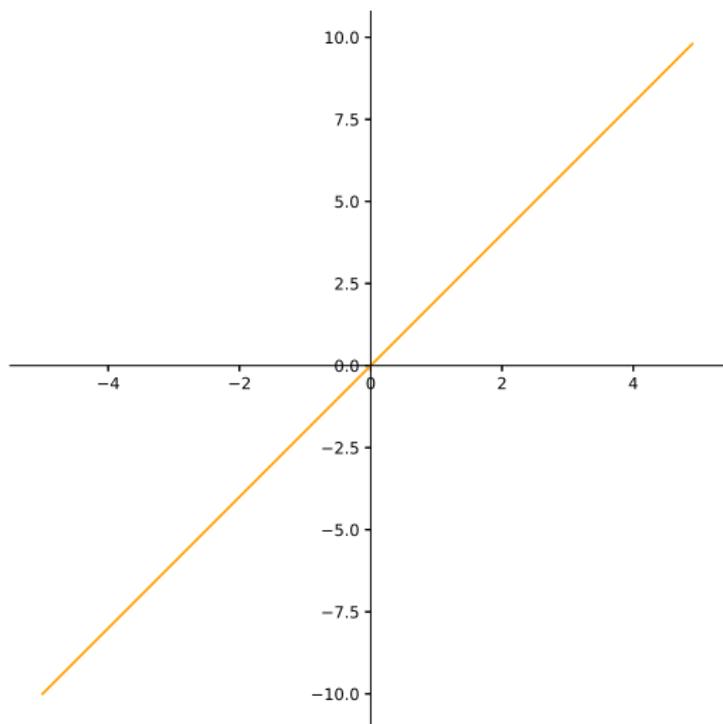
- Porém, a derivada da função valor absoluto é indefinida no ponto 0:



- Por outro lado, a função quadrática apresenta uma única solução ótima, e é bem definida:

$$e(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i \cdot \mathbf{w})^2,$$

Gradiente Descendente



- Temos então que:

$$\nabla e(\mathbf{w}) = -\frac{2}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i \cdot \mathbf{w}) \cdot x_i,$$

- Ou com o mesmo ótimo:

$$\nabla e(\mathbf{w}) = -\frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i \cdot \mathbf{w}) \cdot \mathbf{x}_i,$$

- O novo valor para \mathbf{w} pode ser calculado como:

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \alpha E((y_i - \mathbf{x}_i \cdot \mathbf{w}) \cdot x_i),$$

Gradiente Descendente

- O algoritmo pode ser descrito como:

```
1 gradDesc' :: [[Double]] -> [Double] -> [Double]
2           -> Double -> Int -> [Double]
3 gradDesc' x y w alpha it
4   | convergiu = w
5   | otherwise = gradDesc' x y w' alpha (it-1)
6   where
7     w'          = w .+. (alpha *. nabla)
8     nabla       = mediaVetor
9                 $ map (\(yi, xi) -> yi *. xi)
10                    $ zip (y .-. y') x
11     y'          = map (dotprod w) x
12     convergiu  = (erro' < 1e-6) || (it==0)
```

- Se X é muito grande, temos que paralelizar o cálculo de nabla. Note que o cálculo de y' do qual nabla depende, já está pronto para ser paralelizado.

Gradiente Descendente

- Uma forma de se pensar nessa paralelização é através de um processo mapper e reducer que recebem tuplas de um vetor x e um escalar y :

```
1 gradDesc' :: ChunksOf [(Vector Double, Double)]
2           -> Vector Double -> Double -> Int
3           -> Vector Double
4 gradDesc' xy w alpha it
5   | convergiu = w
6   | otherwise = gradDesc' xy w' alpha (it-1)
7   where
8     w'          = w .+. (alpha *. (nabla ./ n))
9     nabla       = mapReduce mapper reducer xy
10    convergiu   = (it==0)
11    n           = sum $ map length' xy
```

- Escreva as funções mapper e reducer:

Exercício

- Escreva as funções mapper e reducer:

```
1 mapper :: ([Double], Double) -> [Double] -> [Double]
2 mapper (xi, yi) = ( yi - (dotprod w xi) ) *. xi
3
4 reducer :: [Double] -> [Double] -> [Double]
5 reducer = (.+.)
```

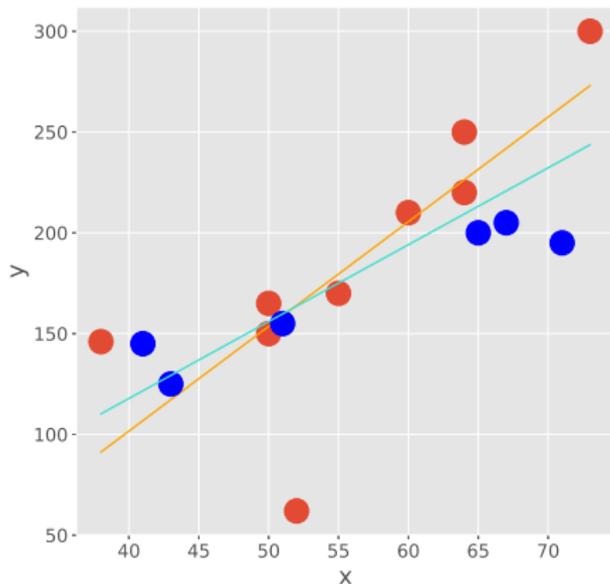
- A escolha do valor de passo de atualização de w é importante pois um valor muito baixo ocasiona uma demora para atingir o ponto de convergência, por outro lado um valor muito alto faz com que o método do gradiente ultrapasse o ponto de ótimo levando a divergência.

Overfitting

- Em muitos casos, a amostra de dados coletada não é representativa. No nosso exemplo, isso pode acontecer ao coletar dados de imóveis de apenas uma região específica.

Overfit

- Isso faz com que a reta de regressão não necessariamente represente a realidade:



- Isso pode ocorrer por causa de:
 - Coleta de dados sem o devido cuidado estatístico.
 - Apenas algumas da d variáveis são realmente importantes.
 - Algumas das variáveis são correlacionadas (ex.: área útil e área total).

- Para resolver tal problema, utilizamos a regularização na função-objetivo. Para isso alteramos o cálculo do erro para:

$$e(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i \cdot \mathbf{w})^2 + \lambda \sum_{i=0}^d |w_i|^p,$$

- Com isso a equação de atualização de \mathbf{w} passa a ser calculada como:

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \alpha E((y_i - \mathbf{x}_i \cdot \mathbf{w}) \cdot x_i) - \lambda l_p,$$

- com l_p sendo o gradiente da regularização escolhida.

- Para $p = 2$ temos a norma quadrática cujo gradiente resultará em:

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \alpha E((y_i - \mathbf{x}_i \cdot \mathbf{w}) \cdot \mathbf{x}_i) - \lambda \mathbf{w}.$$

- Essa é conhecida como regularização l_2 ou ridge. Esse tipo de regularização incentiva baixos valores para todos os elementos do vetor \mathbf{w} .

Regularização l_1

- Para $p = 1$ temos o valor absoluto, o que resulta em descontinuidade. Portanto a fórmula de atualização fica:

$$\Delta \mathbf{w}^t = \mathbf{w}^t + \alpha E((y_i - \mathbf{x}_i \cdot \mathbf{w}) \cdot \mathbf{x}_i)$$
$$\mathbf{w}_i^{t+1} = \begin{cases} \Delta \mathbf{w}_i^t - \lambda & \text{se } \Delta \mathbf{w}_i^t > \lambda \\ \Delta \mathbf{w}_i^t + \lambda & \text{se } \Delta \mathbf{w}_i^t < -\lambda \\ \Delta \mathbf{w}_i^t & \text{c.c.} \end{cases}$$

- Essa é conhecida como regularização l_1 ou lasso. Esse tipo de regularização incentiva que alguns valores de \mathbf{w} sejam zeros.

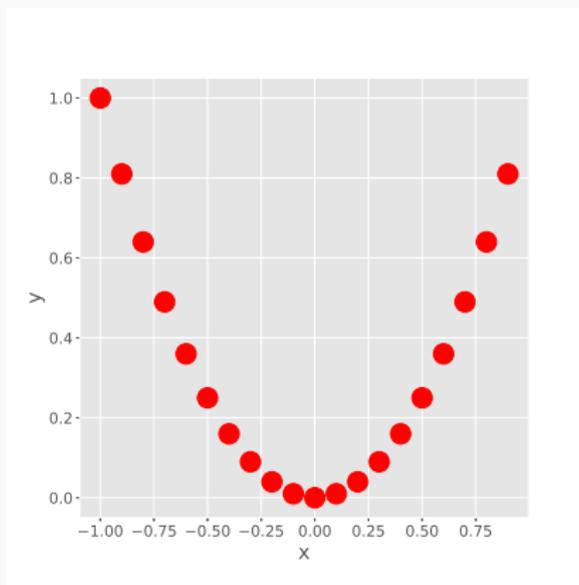
- Podemos combinar ambas as regularizações:

$$\Delta \mathbf{w}^t = \mathbf{w}^t + \alpha E((y_i - \mathbf{x}_i \cdot \mathbf{w}) \cdot x_i)$$
$$\mathbf{w}_i^{t+1} = \begin{cases} \Delta \mathbf{w}_i^t - \lambda_1 - \lambda_2 \mathbf{w} & \text{se } \Delta \mathbf{w}_i^t > \lambda \\ \Delta \mathbf{w}_i^t + \lambda_1 - \lambda_2 \mathbf{w} & \text{se } \Delta \mathbf{w}_i^t < -\lambda \\ \Delta \mathbf{w}_i^t - \lambda_2 \mathbf{w} & \text{c.c.} \end{cases}$$

- Essa é conhecida como regularização ElasticNet.

Variáveis Não-Lineares

Variáveis Não-Lineares



- É fácil perceber que a regressão linear não é um bom modelo para os pontos acima. Qualquer reta utilizada para aproximar os dados irá gerar um erro muito grande em algumas regiões desse gráfico.

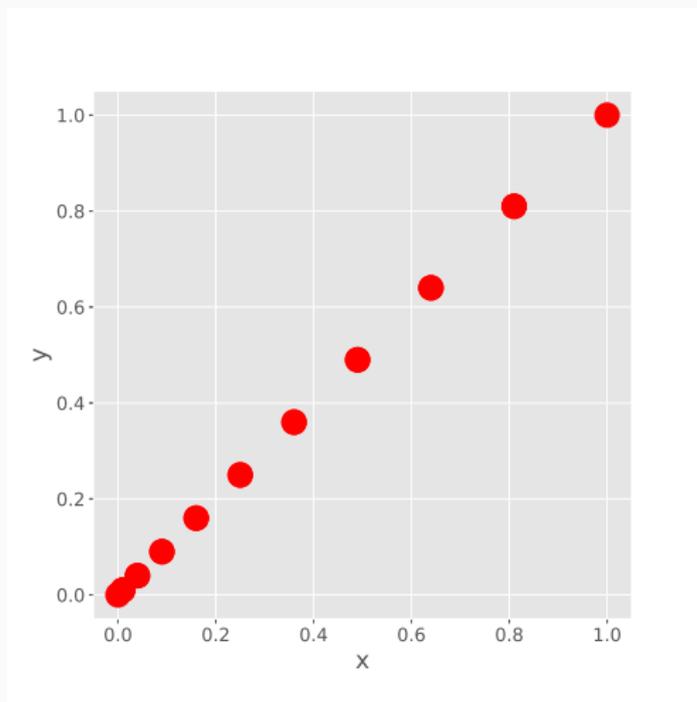
- Uma alternativa para esses casos é criação de novos atributos como funções não lineares.
- Esses novos atributos podem ser:
 - **Transformação:** aplicação de uma função não-linear em uma combinação linear das variáveis.
 - **Interação:** interação polinomial entre duas ou mais variáveis.

- A aplicação de uma função não-linear, se escolhida adequadamente, pode gerar novos atributos que possuem relação linear com a variável alvo.
- Funções comumente utilizadas para essa tarefa:
 - **Sigmoid:** $\frac{1}{1+e^{-x}}$ cria uma variável em um tipo sinal, com valores entre 0 e 1.
 - **Tangente Hiperbólica:** $\tanh(x)$ idem ao anterior, mas variando entre -1 e 1 .
 - **Logaritmo:** $\log(x)$ lineariza variáveis que seguem uma lei de potência.
 - **Rectified Linear Units:** $\max(0, x)$, elimina os valores negativos, relacionado com Redes Neurais.
 - **Softmax:** $\frac{e^{x_i}}{\sum_j e^{x_j}}$, similar a sigmoid, mas também faz com que a soma dos valores de x seja igual a 1.

- A interação polinomial gera interações de grau p entre as variáveis.
- Por exemplo, em um problema com 2 variáveis e $p = 2$ teríamos as variáveis $x_1 \cdot x_1, x_1 \cdot x_2, x_2 \cdot x_2$.
- Já para $p = 3$ teríamos $x_1^3, x_1^2 \cdot x_2, x_1 \cdot x_2^2, x_2^3$.

Interação

- Introduzindo a variável x_1^2 em nosso exemplo, temos:



- k-Vizinhos Mais Próximos
 - Sua adaptação para regressão é trivial, basta fazermos as médias dos valores $\{y_i\}_{i=1}^k$ dos k vizinhos
 - Média aritmética ponderada?

- Levar em consideração todos os pontos, mas ponderando a distância de todos
- Kernel é a função de ponderação
 - Gaussiano é comum: $w(q; x, \sigma) = e^{-\frac{(x-q)^2}{\sigma^2}}$
- Vantagens: kernel permite aproximar funções mais complexas

Kernel Regression: Exemplo

- Considere a base de dados:

x	y
1	1
2	2
3	4
4	8
5	4
6	2
7	1

- Com kernel: $w(q; x) = \frac{1}{(x-q)^2}$
- Como ponto de teste/consulta considere $q = 3.5$.

Kernel Regression: Exemplo

x_i	1	2	3	4	5	6	7
y_i	1	2	4	8	4	2	1
w_i	4/25	4/9	4	4	4/9	4/25	4/49
$w_i \cdot y_i$	4/25	8/9	16	32	16/9	8/25	4/49